

RNA CoSSMos: Characterization of Secondary Structure Motifs—a searchable database of secondary structure motifs in RNA three-dimensional structures

Pamela L. Vanegas, Graham A. Hudson, Amber R. Davis, Shannon C. Kelly, Charles C. Kirkpatrick and Brent M. Znosko*

Department of Chemistry, Saint Louis University, Saint Louis, MO 63103, USA

Received August 22, 2011; Accepted October 12, 2011

ABSTRACT

RNA secondary structure is important for designing therapeutics, understanding protein–RNA binding and predicting tertiary structure of RNA. Several databases and downloadable programs exist that specialize in the three-dimensional (3D) structure of RNA, but none focus specifically on secondary structural motifs such as internal, bulge and hairpin loops. The RNA Characterization of Secondary Structure Motifs (RNA CoSSMos) database is a freely accessible and searchable online database and website of 3D characteristics of secondary structure motifs. To create the RNA CoSSMos database, 2156 Protein Data Bank (PDB) files were searched for internal, bulge and hairpin loops, and each loop's structural information, including sugar pucker, glycosidic linkage, hydrogen bonding patterns and stacking interactions, was included in the database. False positives were defined, identified and reclassified or omitted from the database to ensure the most accurate results possible. Users can search via general PDB information, experimental parameters, sequence and specific motif and by specific structural parameters in the subquery page after the initial search. Returned results for each search can be viewed individually or a complete set can be downloaded into a spreadsheet to allow for easy comparison. The RNA CoSSMos database is automatically updated weekly and is available at <http://cossmos.slu.edu>.

INTRODUCTION

Deoxyribonucleic acid (DNA) and ribonucleic acid (RNA) are fundamental components of all cellular life, and much of their functionality is due to their 3D structure. While DNA is double stranded, RNA is normally single stranded and folds back upon itself, creating sites of mismatched nucleotides within the center and at the end of a duplex. These loops, including internal, bulge and hairpin loops, contribute to the functionality as well as the overall tertiary folding of the nucleic acid. Knowledge of these specific structural motifs can aid in the design of therapeutics (1,2), the understanding of RNA–protein interactions (3,4) and the prediction of tertiary structure folding of RNA (5,6). For example, myotonic dystrophy types 1 and 2 (DM1 and DM2) are diseases occurring from the expression of a tandem repeat in the genome resulting in a 1×1 U•U or a 2×2 $5'$ CU $3'/3'$ UC $5'$ symmetric internal loop, respectively. Specifically-designed small molecules could target these loops and act as effective therapeutics against these two diseases (1). Additionally, C•A 1×1 internal loops are a recognition site for adenosine deaminase acting on RNA type 2 (ADAR2) (7) and can be targeted with small molecules to help prevent misediting by the enzyme (2). In some RNA macromolecules, bulge loops and hairpin loops play an important role in protein–RNA binding (3,4), such as the $5'$ CCCG $3'$ hairpin in pre-mRNA that is selectively bound by nucleolin (8).

Related databases and programs

Currently, there are several databases that focus on RNA 3D structure. RNAFRABase 2.0 (9,10) is a powerful online tool capable of analyzing large oligomers and

*To whom correspondence should be addressed. Tel: +1 314 977 8567 ; Fax: +1 314 977 2521; Email: znoskob@slu.edu

displaying the torsional angles and coordinates. RNAFRABase also allows users to input structural information and displays selected secondary structures. RNAJunction (11) contains structural information about three way junctions and kissing hairpins found in the Protein Data Bank (PDB) (12), while SCOR (13) categorizes several hundred PDB files by both structural and functional classifications. FR3D (14) and the online version, WebFR3D (15), give users the ability to specify sequence and structural characteristics of desired RNA 3D motifs in order to run a real-time scan of PDB files. Additionally, several programs use sequence alignment against a known 3D structure in order to predict an unknown tertiary structure (16–23). Others scan through the 3D structure of a PDB file to identify local motifs within the specific file (24–28).

The RNA Characterization of Secondary Structure Motifs (RNA CoSSMos) database can be used by researchers as a complementary resource to those previously mentioned. The RNA CoSSMos database focuses on the 3D characteristics of secondary structure motifs in RNA; these include symmetric and asymmetric internal loops, bulge loops and hairpin loops. The structural information is pre-compiled to allow for faster searches, and the graphical user interface is designed to be intuitive, limiting pre-requisite knowledge of specialized syntax. False positives, as defined later, have been omitted from the database or reclassified to create the most accurate database possible. RNA CoSSMos is automatically updated with new PDB files and their structural characterizations weekly, keeping the database current.

DATABASE DESIGN AND CONTENT

Data extraction and motif identification

The extraction of the 3D RNA structures from the PDB was similar to that previously described by Davis *et al.* (29). Briefly stated, all PDB structures containing RNA, including RNA–protein structures and RNA–DNA hybrids, were downloaded. At the time of publication, 2156 PDB structures were included. Input descriptors were written based upon the 1×1 symmetric internal mismatch input descriptor designed by Davis *et al.* (29). An input descriptor is necessary to designate which base pairings are valid, e.g. canonical base pairs for nearest neighbors and non-canonical base pairs for the mismatch, and which base pairings are invalid, e.g. non-canonical base pairs for nearest neighbors and canonical base pairs for the mismatch. *MC-Search* (24,25) was used to identify motifs based upon these input descriptors and then to ‘clip’ the motif so all nucleotides that were not directly part of the loop or the closing base pair(s) were removed. *MC-Search* works by searching through the three-dimensional structure for specific nucleotide interactions, defined by Saenger (30), Westhof (31,32) or Major (25). After *MC-Search* ‘clipped’ the structure, *MC-Annotate* 1.6.2 (24,25) was then used to determine structural data including sugar puckers, base pairings and stacking interactions between the closing base pairs and the loop nucleotides. Standard code was used to

process *MC-Annotate* output and compile it into the database.

Removal of false positives found by *MC-Search*

A unique feature of the RNA CoSSMos database is the removal of false positives found by *MC-Search* before the data is exported. For the purposes of the RNA CoSSMos database, there are two criteria for a false positive; if either is met, then the data are excluded from that subcategory of the database and reclassified. The first criterion requires that the nearest neighbors of the loop must be an A–U, C–G, G–C, G–U, U–A or U–G base pair. The second criterion requires that the mismatched nucleotides themselves cannot potentially form an A–U, C–G, G–C, G–U, U–A, or U–G base pair. These criteria and definitions of false positives appear to be unique to the RNA CoSSMos database; therefore, for researchers who agree with these criteria, the reclassification of false positives offers a distinct advantage over other available databases and software.

In symmetric internal loops, for instance, *MC-Search* may classify $\begin{bmatrix} 5' \text{GAAAC} 3' \\ 3' \text{CAAUG} 5' \end{bmatrix}$ as a 3×3 symmetric internal loop, with the mismatch nucleotides underlined. According to the previously stated definition, this would be considered a false positive on the basis that the last mismatch of the loop, the 3' adenine of the upper strand and the 5' uracil of the lower strand, has the ability to form an A–U pair. In the RNA CoSSMos database, this has been reclassified as the 2×2 symmetric internal loop $\begin{bmatrix} 5' \text{GAAAC} 3' \\ 3' \text{CAAU} 5' \end{bmatrix}$. As shown here, symmetric internal loops that are determined to be false positives are reclassified as either smaller symmetric loops or as having no mismatches.

False positives are also possible for asymmetric internal loops. With asymmetric internal loops, either smaller asymmetric internal loops or bulges are possible results from the screenings. For example, the sequence $\begin{bmatrix} 5' \text{CACA} 3' \\ 3' \text{GAAGU} 5' \end{bmatrix}$ was classified as a 2×3 internal loop by *MC-Search*. Using the same methodology as with the symmetric internal loops described above, this would be classified in the RNA CoSSMos database as the 1×2 internal loop $\begin{bmatrix} 5' \text{CAC} 3' \\ 3' \text{GAAG} 5' \end{bmatrix}$. Another internal loop classified as a 2×3 internal loop by *MC-Search*, $\begin{bmatrix} 5' \text{AGCC} 3' \\ 3' \text{UCGAG} 5' \end{bmatrix}$, is reclassified in the RNA CoSSMos database as the single nucleotide bulge loop $\begin{bmatrix} 5' \text{CC} 3' \\ 3' \text{GAG} 5' \end{bmatrix}$.

While the majority of the asymmetric false positives are able to be categorized into a specific smaller loop or bulge, some remained that resulted in an indeterminate mismatch. In order to maintain equivalent data analysis between both the symmetric and asymmetric internal loops, the undetermined loops were excluded from the RNA CoSSMos database. For instance, the 2×3 internal loop classified by *MC-Search* $\begin{bmatrix} 5' \text{AUAC} 3' \\ 3' \text{UAUUG} 5' \end{bmatrix}$ is a false positive; however, this loop was excluded from the

RNA CoSSMos database. Due to the possibility of the adenosine in the top strand forming a base pair with either of the uracils on the bottom strand, which would result in two distinct single nucleotide bulges, loops such as this would need to be considered on an individual basis and therefore were omitted from the RNA CoSSMos database. Hairpin loops of 5–7 nt were evaluated for false positives in a similar manner. Bulge loops were assessed to ensure the closing base pairs were an A–U, C–G, G–C, G–U, U–A or U–G base pair. Within the RNA CoSSMos database after reclassifications, there are 11,860 symmetric internal loops, 7,781 asymmetric internal loops, 30,231 bulge loops, and 22,257 hairpin loops (Table 1).

Database search capabilities and results output

There are several ways of searching the RNA CoSSMos database. Fields can be used in any combination and include general PDB information, experimental parameters, motif and sequence (Figure 1A). General PDB information parameters include PDB identification number, the authors of the published structure and keywords found within the PDB file. It is also possible to search by the type of experiment that was used to determine the 3D structure; X-ray diffraction and cryo-electron microscopy can be limited by the resolution of the experiment, while NMR experiments can be limited by the number of structures within the ensemble. Additionally,

after the initial search has been run, it is possible to specify preferred structural characteristics in the subquery.

Along with these optional experimental parameter filters, users have the ability to search for specific motifs in differing sizes. These motifs include 1×1 , 2×2 , 3×3 , 4×4 and 5×5 symmetric internal loops, 1×2 , 1×3 , 1×4 , 1×5 , 2×3 , 2×4 , 2×5 , 3×4 , 3×5 and 4×5 asymmetric internal loops, hairpin loops of 3, 4, 5, 6 or 7 nt and bulge loops of 1, 2, 3, 4 or 5 nt. Within the motif-specific search, it is possible to search for one submotif, several submotifs or all submotifs. For example, selecting only triloops, both triloops and tetraloops, or all hairpins found within the database is allowed by RNA CoSSMos. To search more precisely, the desired sequence can be selected by distinguishing closing base pairs and mismatched nucleotides using the seven standard base abbreviations: *A*, *C*, *G*, *U*, *R* (any purine), *Y* (any pyrimidine) and *N* (any nucleotide). Once the results from the initial search have been returned, the user can choose to view those results (Figure 1B) or modify the parameters (Figure 2B). Search modifications can be made by either narrowing the standard search parameters or by using subqueries, which allow the user to search through the previous dataset and specify certain structural features, i.e. sugar puckers, glycosidic linkages, interacting edges and stacking conformations (Figure 2B). The subquery page can be located through a tab at the top of the results page.

Within the results pages of the RNA CoSSMos database, the user can look at each specific mismatch found in the PDB or at the entire dataset via the download results option, which exports as a pound delimited file into a spreadsheet. The individual detailed results are useful for specific structural searches, while the downloadable results allow the user to more easily compare groups of structures. As explained in Davis *et al.* (29), *MC-Annotate* characterizes each mismatch with four different parameters: residue conformation, base pairing, adjacent stacking and non-adjacent stacking. For each nucleotide, the residue conformation is given by describing the sugar pucker as either *endo* or *exo* and the glycosidic linkage as either *syn* or *anti*. For the base pairing characterizations that are labeled as ‘Interacting Edges’ in the RNA CoSSMos database, *MC-Annotate* relies on the nomenclature schemes developed by Major (25), Saenger (30) and Westhof (31,32). Abbreviations [i.e. Watson–Crick (W), Hoogsteen (H) and sugar (S)] are used to identify the edge of the base that is involved in the hydrogen bonding. In most cases, a two-letter designation is used to describe the interacting edges where the first letter describes which edge of the base is pairing, and the second specifies where on that edge the bonding is occurring (25). For example, *Ws* would be used to designate that the Watson–Crick side of the base near the sugar edge is participating in the hydrogen bond. Some bonding patterns are also designated by either a Roman numeral (Saenger notation) (30) or an Arabic numeral (Westhof notation) (31,32). For interacting edges that are not designated with the two-letter abbreviation, a different, non-standard annotation may be used to describe the hydrogen-bonding patterns, such as

Table 1. Number of motifs found in the RNA CoSSMos database

Loop Type and Size	Number of Loops in RNA CoSSMos
Internal Symmetric Loops	
1×1	7 545
2×2	1 993
3×3	1 568
4×4	558
5×5	196
Internal Asymmetric Loops	
1×2	3 329
1×3	817
1×4	553
1×5	55
2×3	1 647
2×4	417
2×5	133
3×4	372
3×5	69
4×5	389
Hairpin Loops	
3	3 358
4	10 225
5	4 498
6	2 746
7	1 430
Bulge Loops	
1	21 070
2	5 941
3	2 025
4	489
5	706



Figure 1. The CoSSMos (A) search page and (B) results page.

$5' \text{O}2\text{P}/3' \text{Bh}$ to describe a bifurcated hydrogen-bonding pattern between the $2'$ oxygen of the sugar and the Hoogsteen face. Stacking interactions, both adjacent and non-adjacent, are described by the terms upward, downward, outward or inward as proposed by Major and Thibault (33). All mismatches within the database are characterized by these four designations.

The detailed results page (Figure 2A) for each motif displays all structural characterizations for every applicable nucleotide and a 'clipped' PDB structure in a Jmol (34) applet, which allows for any user with a web browser that supports Java to view the structure. In order to ensure the accurate representation of the motif structure, the PDB file for the 'clipped' structure was not altered;

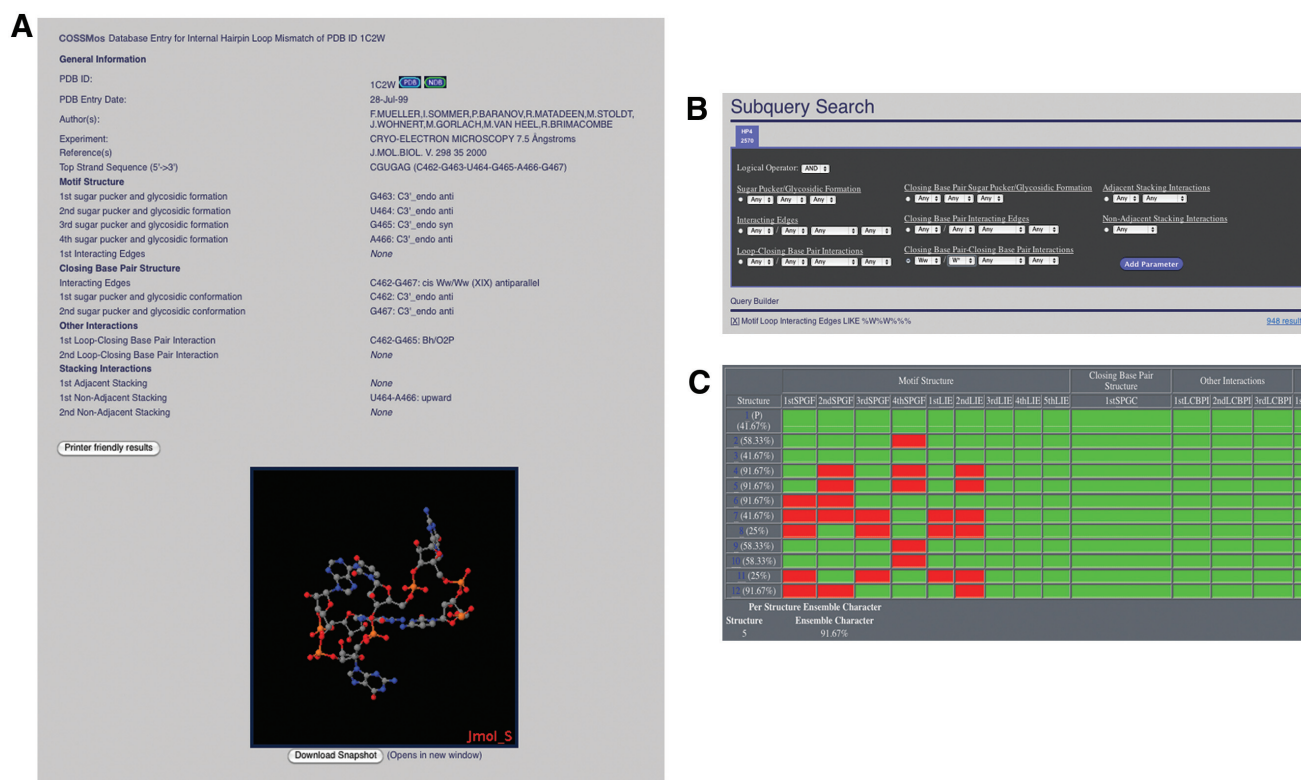


Figure 2. The CoSSMos (A) detailed results page, (B) subquery search page, and (C) the NMR overview page.

instead, the nucleotides and amino acids that are not involved in the mismatch are blacked out, reducing the possibility of unintentionally deleting or modifying the structural data. This ‘clipped’ structure can be manipulated as in any Jmol program and then downloaded as a snapshot. In addition to the structural information and the clipped structure in Jmol, the detailed results page includes an option for printer-friendly results, links to the PDB (12) and links to the Nucleic Acid Database (35), where applicable.

Additional features

Beyond the search capabilities and the results output, the RNA CoSSMos database contains an optional username, a frequently asked questions page and an NMR overview page. Although the RNA CoSSMos database is freely accessible to all, registering with a username is encouraged. Upon registering, users will be able to save up to 10 searches on the RNA CoSSMos database, as well as receive emails about important updates to the RNA CoSSMos database. The frequently asked questions page is linked to the home page of the website and contains information on both the RNA CoSSMos database itself and the structural characterizations, including definitions of and graphics depicting the interactions between nucleotides. For all ensemble structures determined by NMR, the RNA CoSSMos database contains a unique NMR overview page, providing a convenient way to analyze the structural characteristics between ensemble structures (Figure 2C). In the NMR

overview, the first structure of the ensemble is arbitrarily defined as the prime structure to which the others are compared. The remaining structures’ characteristics are evaluated and color coded based upon the equivalency. If the interaction is equivalent to that in the prime structure, the box is colored green; if not, the box is colored red. Additionally, each structure is assigned a percentage based upon its contribution to the ensemble character. The structure with the highest percentage may be considered as the representative structure of the ensemble.

CONCLUSIONS AND FUTURE DIRECTIONS

The RNA CoSSMos database is a unique online tool that gives researchers the ability to search for 3D characteristics of RNA secondary structure motifs without creating the need for the user to run a comprehensive search. The database design is a simple and intuitive graphical user interface, which eliminates the need for complex syntax. Within the database, the reclassification of false positives eliminates incorrect identifications of the motifs, allowing for the most accurate database possible. Additionally, multiple search parameters, including sequence and motif structure, make RNA CoSSMos versatile for many different uses. The downloadable results, the detailed results pages and the NMR overview pages create many different options for viewing the structural characterizations of the motifs. Future versions of RNA CoSSMos may extend into DNA secondary structure motifs and higher order RNA motifs,

such as A-platforms and U-turns. Additionally, the indeterminate mismatches may be included, and users will be able to search for them. Future directions of the database will also be driven by feedback from the users. As the PDB continues to grow, so will the RNA CoSSMos database and its capabilities.

AVAILABILITY

The RNA CoSSMos database is freely available online at <http://cossmos.slu.edu>. Users of the RNA CoSSMos database should cite this article and are encouraged to cite the original references for *MC-Search* and *MC-Annotate* from Francois Major's laboratory (24,25).

ACKNOWLEDGEMENTS

The authors would like to thank Francois Major for providing us with executable versions of *MC-Search* and *MC-Annotate* and for providing assistance with the software in the context of the work described here. Additionally, the authors would like to thank Nina Zulic Hausmann for her initial work on the input descriptors.

FUNDING

The National Institutes of Health (1R15GM085699-01A1 to B.M.Z.). Funding for open access charge: NIH.

Conflict of interest statement. None declared.

REFERENCES

- Lee, M.M., Childs-Disney, J.L., Pushechnikov, A., French, J.M., Sobczak, K., Thornton, C.A. and Disney, M.D. (2009) Controlling the specificity of modularly assembled small molecules for RNA via ligand module spacing: targeting the RNAs that cause myotonic muscular dystrophy. *J. Am. Chem. Soc.*, **131**, 17464–17472.
- Tran, T. and Disney, M.D. (2011) Molecular recognition of 6'-N-5-hexynoate kanamycin A and RNA 1x1 internal loops containing CA mismatches. *Biochemistry*, **50**, 962–969.
- Gupta, A. and Gribskov, M. (2011) The role of RNA sequence and structure in RNA-protein interactions. *J. Mol. Biol.*, **409**, 574–587.
- Hermann, T. and Patel, D.J. (2000) RNA bulges as architectural and recognition motifs. *Structure*, **8**, 47–54.
- Leontis, N.B. and Westhof, E. (2003) Analysis of RNA motifs. *Curr. Opin. Struct. Biol.*, **13**, 300–308.
- Reiter, N.J., Chan, C.W. and Mondragón, A. (2011) Emerging structural themes in large RNA molecules. *Curr. Opin. Struct. Biol.*, **21**, 319–326.
- Wong, S.K., Sato, S. and Lazinski, D.W. (2001) Substrate recognition by ADAR1 and ADAR2. *RNA*, **7**, 846–858.
- Allain, F.H.T., Bouvet, P., Dieckmann, T. and Feigon, J. (2000) Molecular basis of sequence-specific recognition of pre-ribosomal RNA by nucleolin. *EMBO J.*, **19**, 6870–6881.
- Popenda, M., Blazewicz, M., Szachniuk, M. and Adamiak, R.W. (2008) RNA FRABASE version 1.0: an engine with a database to search for the three-dimensional fragments within RNA structures. *Nucleic Acids Res.*, **36**, D386–D391.
- Popenda, M., Szachniuk, M., Blazewicz, M., Wasik, S., Burke, E.K., Blazewicz, J. and Adamiak, R.W. (2010) RNA FRABASE 2.0: an advanced web-accessible database with the capacity to search the three-dimensional fragments within RNA structures. *BMC Bioinformatics*, **11**, 231–242.
- Bindewald, E., Hayes, R., Yingling, Y., Kasprzak, W. and Shapiro, B.A. (2008) RNAJunction: a database of RNA junctions and kissing loops for three-dimensional structural analysis and nanodesign. *Nucleic Acids Res.*, **36**, D392–D397.
- Berman, H.M., Westbrook, J., Feng, Z., Gilliland, G., Bhat, T.N., Weissig, H., Shindyalov, I.N. and Bourne, P.E. (2000) The Protein Data Bank. *Nucleic Acids Res.*, **28**, 235–242.
- Klosterman, P.S., Tamura, M., Holbrook, S.R. and Brenner, S.E. (2002) SCOR: a structural classification of RNA database. *Nucleic Acids Res.*, **30**, 392–394.
- Sarver, M., Zirbel, C.L., Stombaugh, J., Mokdad, A. and Leontis, N.B. (2008) FR3D: finding local and composite recurrent structural motifs in RNA 3D structures. *J. Math. Biol.*, **56**, 215–252.
- Petrov, A.I., Zirbel, C.L. and Leontis, N.B. (2011) WebFR3D—a server for finding, aligning and analyzing recurrent RNA 3D motifs. *Nucleic Acids Res.*, **39**, W50–W55.
- Dror, O., Nussinov, R. and Wolfson, H. (2005) ARTS: alignment of RNA tertiary structures. *Bioinformatics*, **21**, ii47–ii53.
- Dror, O., Nussinov, R. and Wolfson, H.J. (2006) The ARTS web server for aligning RNA tertiary structures. *Nucleic Acids Res.*, **34**, W412–W415.
- Ferrè, F., Ponty, Y., Lorenz, W.A. and Clote, P. (2007) DIAL: a web server for the pairwise alignment of two RNA three-dimensional structures using nucleotide, dihedral angle and base-pairing similarities. *Nucleic Acids Res.*, **35**, W659–W668.
- Chang, Y.-F., Huang, Y.-L. and Lu, C.L. (2008) SARSA: a web tool for structural alignment of RNA using a structural alphabet. *Nucleic Acids Res.*, **36**, W19–W24.
- Capriotti, E. and Marti-Renom, M.A. (2008) RNA structure alignment by a unit-vector approach. *Bioinformatics*, **24**, i112–i118.
- Lai, C.-E., Tsai, M.-Y., Liu, Y.-C., Wang, C.-W., Chen, K.-T. and Lu, C.L. (2009) FASTR3D: a fast and accurate search tool for similar RNA 3D structures. *Nucleic Acids Res.*, **37**, W287–W295.
- Rother, M., Rother, K., Puton, T. and Bujnicki, J.M. (2011) ModeRNA: a tool for comparative modeling of RNA 3D structure. *Nucleic Acids Res.*, **39**, 4007–4022.
- Kirillova, S., Tosatto, S.C.E. and Carugo, O. (2010) Open access FRASS: the web-server for RNA structural comparison software. *BMC Bioinformatics*, **11**, 327–334.
- Gendron, P., Lemieux, S. and Major, F. (2001) Quantitative analysis of nucleic acid three-dimensional structures. *J. Mol. Biol.*, **308**, 919–936.
- Lemieux, S. and Major, F. (2002) RNA canonical and non-canonical base pairing types: a recognition method and complete repertoire. *Nucleic Acids Res.*, **30**, 4250–4263.
- Duarte, C.M. (2003) RNA structure comparison, motif search and discovery using a reduced representation of RNA conformational space. *Nucleic Acids Res.*, **31**, 4755–4761.
- Zhong, C., Tang, H. and Zhang, S. (2010) RNAMotifScan: automatic identification of RNA structural motifs using secondary structural alignment. *Nucleic Acids Res.*, **38**, 176e.
- Macke, T.J., Ecker, D.J., Gutell, R.R., Gautheret, D., Case, D.A. and Sampath, R. (2001) RNAMotif, an RNA secondary structure definition and search algorithm. *Nucleic Acids Res.*, **29**, 4724–4735.
- Davis, A.R., Kirkpatrick, C.C. and Znosko, B.M. (2010) Structural characterization of naturally occurring RNA single mismatches. *Nucleic Acids Res.*, **39**, 1081–1094.
- Saenger, W. (1984) *Principles of Nucleic Acid Structure*. Springer, New York.
- Leontis, N.B. and Westhof, E. (1998) Conserved geometrical base-pairing patterns in RNA. *Q. Rev. Biophys.*, **31**, 399–455.
- Leontis, N.B. and Westhof, E. (2001) Geometric nomenclature and classification of RNA base pairs. *RNA*, **7**, 499–512.
- Major, F. and Thibault, P. (2007) In: Lengauer, T. (ed.), *Bioinformatics: From Genetics to Therapies*. Wiley-VCH, Weinheim, Germany, pp. 491–539.
- Jmol: an open source Java viewer for chemical structures in 3D. <http://www.jmol.org/>.
- Berman, H.M., Olson, W.K., Beveridge, D.L., Westbrook, J., Gelbin, A., Demeny, T., Hsieh, S.H., Srinivasan, A.R. and Schneider, B. (1992) The Nucleic Acid Database: a comprehensive relational database of three-dimensional structures of nucleic acids. *Biophys. J.*, **63**, 751–759.